



# Spectral sensitivity of products of bidiagonals <sup>1</sup>

Beresford N. Parlett <sup>2</sup>

*Mathematics Department and Computer Science Division, EECS Department, University of California, Berkeley, CA 94720, USA*

Received 28 November 1996; accepted 18 August 1997

Submitted by V. Mehrmann

---

## Abstract

Simple formulae are derived for certain relative condition numbers. Tridiagonal matrices may be represented as products of bidiagonals in various ways depending on properties such as symmetry and positive definiteness. The condition numbers give the amplification factor for relative changes in a nonzero eigenvalue caused by relative changes in an entry of a bidiagonal factor. The formulae show that in many, but not all cases these condition numbers are of modest size. Several examples illustrate the results and raise new questions. © 1998 Elsevier Science Inc. All rights reserved.

---

## 1. Introduction

Any engineer who is aware of the limitations on his data must be interested in how well the output of some process is determined by the data. In this study the output is a set of eigenvalues and the input is usually a square matrix. Perturbation theory teaches us some hard lessons. In general some eigenvalues may be exceedingly sensitive to small changes in the matrix entries while others may be the opposite. For Hermitian matrices, and these include the real symmetric matrices, no eigenvalue can change by more than the norm of the change in the matrix (H. Weyl). See Chapter 10 of [1] for a simple proof. However, in some applications, the spectrum ranges over many orders of magnitude and it is usually the eigenvalues near 0 that are of interest. Unfortunately

---

<sup>1</sup> This research paper is based on a talk given at the minisymposium on Perturbation Theory during the 6th ILAS Conference in Chemnitz, Germany, in August 1996.

<sup>2</sup> Supported by ONR, Contract N000014-90-J-1372.

Weyl's theorem does not guarantee that these small eigenvalues have as many correct decimal digits as the large (unwanted) ones.

In a number of important cases a generalized Hermitian problem  $(K - \lambda M)\mathbf{u} = \mathbf{0}$  does determine its small eigenvalues to high relative accuracy but this valuable property has been lost by heeding the mathematician's advice to reduce a difficult problem (say the pair  $K, M$ ) to an easier one (a standard eigenproblem  $A - \lambda I$  with  $A = KM^{-1}$  or  $A = M^{-1/2}KM^{-1/2}$ ). There is more on this topic in Section 4.

Even symmetric tridiagonal matrices that are positive definite do not always determine their tiny eigenvalues to an acceptable number of figures.

Armed with this sobering knowledge numerical analysts are careful not to try to force their algorithms to achieve unattainable accuracies.

The objective that is repeated, like a mantra, throughout the community of eigenvalue hunters is to achieve 'as much accuracy as the data warrants'.

We acknowledge clever recent research that overlaps the results given here. See [2–5], for example. The formulae in Theorem 1 are new but the 'result' is known, see [6]. Theorem 2 is not as general as results given in [5], but Theorem 3 does break new ground. We show that a little calculus gives a lot of insight. Kahan's ingenious proof, in 1966 (!) but see [6], that all singular values are determined to high relative accuracy by the entries in a bidiagonal matrix was not really necessary! We are not suggesting that traditional perturbation theory be replaced by calculus but people do remember simple approaches and we show how naturally a representation as a product of matrices can lead to expressions for relative accuracy.

Before starting let us note that a single matrix  $M$  may be regarded as a representative of the equivalence class of all its translates  $\{M + \sigma I, \sigma \in \mathbb{C}\}$ . In this sense there is no natural origin for the standard eigenvalue problem. However it is not so easy to represent in factored form the translates of a product  $FG$  of invertible matrices.

## 2. Background

The adjectives absolute and relative will occur frequently in the investigation and it seems advisable to say what is meant.

Let  $B$  be a square matrix that depends smoothly on some parameter  $\tau$ . Following Newton's notation for derivatives let  $\dot{B}$  denote the derivative of  $B$  with respect to  $\tau$  and consider a typical simple eigenvalue  $\lambda$  with column and row eigenvectors  $\mathbf{x}$  and  $\mathbf{y}^*$ . Thus

$$B\mathbf{x} = \lambda\mathbf{x}, \quad \mathbf{y}^*B = \lambda\mathbf{y}^*,$$

and without loss of generality one may assume  $\mathbf{y}^*\mathbf{x} > 0$ . Standard differentiation yields

$$\dot{\lambda} = \frac{\mathbf{y}^* \dot{B} \mathbf{x}}{\mathbf{y}^* \mathbf{x}} = \frac{\|\mathbf{y}^*\| \cdot \|\mathbf{x}\|}{\mathbf{y}^* \mathbf{x}} \cdot \frac{\mathbf{y}^* \dot{B} \mathbf{x}}{\|\mathbf{y}^*\| \cdot \|\mathbf{x}\|}$$

and  $\|\mathbf{y}^*\| \cdot \|\mathbf{x}\| / \mathbf{y}^* \mathbf{x}$  is the (absolute) condition number of  $\lambda$ ;  $\text{cond}(\lambda)$ . In particular the sensitivity of  $\lambda$  to a change in  $b_{jk}$ , i.e.  $\tau = b_{jk}$ , is given by

$$\frac{\partial \lambda}{\partial b_{jk}} = \frac{\mathbf{y}^* \mathbf{e}_j \cdot \mathbf{e}_k^* \mathbf{x}}{\mathbf{y}^* \mathbf{x}} = \frac{\bar{y}(j)x(k)}{\mathbf{y}^* \mathbf{x}},$$

where  $I = (\mathbf{e}_1, \dots, \mathbf{e}_n)$ .

We call  $|\partial \lambda / \partial b_{jk}|$  the *absolute sensitivity* (of  $\lambda$  to  $b_{jk}$ ). We may also be interested in the *relative* sensitivity of  $\lambda$ , for  $\lambda \neq 0$ , i. e.  $|(\partial \lambda / \partial b_{jk}) / \lambda|$ . Even more relevant is the relative change in  $\lambda$  due to small relative changes in the  $(j, k)$  entry, i.e.

$$\left| \frac{\partial \lambda}{\partial b_{jk}} \cdot \frac{b_{jk}}{\lambda} \right|$$

and this is what we mean by the *relative sensitivity*. If  $|(\partial \lambda / \partial b_{jk}) \cdot (b_{jk} / \lambda)| = 10^2$  we would say, somewhat loosely, that uncertainty in the 5th decimal digit of  $b_{jk}$  provokes uncertainty in the 3rd decimal digit of  $\lambda$ . The quantity  $\text{cond}(\lambda)$  defined above is absolute and it is an upper bound on the sensitivity of  $\lambda$  to (absolute) changes in any entry. As mentioned in the introduction, there are examples of positive definite symmetric tridiagonal matrices such that *relative* changes to some off diagonal entries of order  $\epsilon$  produce *relative* changes in some eigenvalues of order  $\sqrt{\epsilon}$ . Thus we cannot be sure, in general, that tiny  $\lambda$  are determined to any correct figures when matrix entries are uncertain or noisy. This is sad but true and prompts us to look for alternative representations of matrices that might determine all eigenvalues well, regardless of magnitude.

Recent investigations have shown that certain classes of matrices do determine their small eigenvalues to high relative accuracy. See [6–10].

### 3. Singular values of bidiagonal

In 1966, W. Kahan proved a surprising result which is in stark contrast to the high relative sensitivities of small eigenvalues to small changes in matrix entries. He showed that small relative changes in the entries of a bidiagonal matrix lead to small relative changes in all the singular values, see [6]. Let

$$B = \text{bidiag} \begin{pmatrix} & b_1 & & & & & \\ a_1 & & b_2 & & & & \\ & a_2 & & \ddots & & & \\ & & & & b_{n-2} & & \\ & & & & & b_{n-1} & \\ & & & & & & a_n \end{pmatrix}$$

have  $(\sigma, \mathbf{v}, \mathbf{u})$  as a typical singular triple:

$$Bv = u\sigma, \quad B^t u = v\sigma, \quad (1)$$

$$v^t v = u^t u = 1. \quad (2)$$

Suppose that  $a_i > 0$ ,  $i = 1, 2, \dots, n$ , and  $b_j > 0$ ,  $j = 1, \dots, n-1$ .

To show that eigenvalues have not been abandoned note that if

$$\frac{\partial \sigma}{\partial a_k} \cdot \frac{a_k}{\sigma} = \kappa,$$

then this is equivalent to a result concerning  $\lambda = \sigma^2$  of  $B^t B$  since

$$\frac{\partial \lambda}{\partial a_k} \cdot \frac{a_k}{\lambda} = 2\sigma \frac{\partial \sigma}{\partial a_k} \cdot \frac{a_k}{\sigma^2} = 2\kappa.$$

Kahan's original proof was not published and uses Sylvester's Inertia theorem in an ingenious way. The proof in [6] is a little more general but follows the same far from obvious reasoning. A simple proof is given in [11].

In contrast our Theorem 1 is straightforward and yields stronger results than Theorem 2 in [6] for tiny perturbations. In our terminology, Kahan bounded the relative condition number by 1 but that bound holds for all relative perturbations both large and small while our Theorem 1 applies only to infinitesimal changes.

In what follows it is convenient to define

$$b_0 = b_n = 0.$$

**Theorem 1.** *With the notation given in Eqs. (1) and (2), since  $\sigma \neq 0$ ,*

$$(a) \quad \frac{\partial \sigma}{\partial a_k} \cdot \frac{a_k}{\sigma} = \sum_{i=1}^k v(i)^2 - \sum_{j=1}^{k-1} u(j)^2 = \sum_{m=k}^n u(m)^2 - \sum_{l=k+1}^n v(l)^2,$$

$$(b) \quad \frac{\partial \sigma}{\partial b_k} \cdot \frac{b_k}{\sigma} = \sum_{i=1}^k (v(i)^2 - u(i)^2) = \sum_{m=k+1}^n (u(m)^2 - v(m)^2).$$

**Proof.** First derive well known expressions for  $\partial \sigma / \partial a_k$  and  $\partial \sigma / \partial b_k$ . From Eqs. (1) and (2),

$$\sigma = u^t B v. \quad (3)$$

Write out Eq. (3) using  $B$ 's entries to find

$$\frac{\partial \sigma}{\partial a_j} = u(j)v(j), \quad \frac{\partial \sigma}{\partial b_j} = u(j)v(j+1). \quad (4)$$

These expressions are simple but do not reveal the dependence on  $\sigma$ ,  $a_j$ , and  $b_j$ . To remedy the situation write out each equation in (1) in detail

$$a_j v(j) + b_j v(j+1) = u(j)\sigma, \quad b_{j-1} u(j-1) + a_j u(j) = v(j)\sigma. \quad (5)$$

Multiply Eq. (4) by  $a_j$  and  $b_j$  and substitute into Eq. (5) to find, for  $j = 1, 2, \dots, n$ ,

$$a_j \frac{\partial \sigma}{\partial a_j} + b_j \frac{\partial \sigma}{\partial b_j} = u(j)^2 \sigma, \quad b_{j-1} \frac{\partial \sigma}{\partial b_{j-1}} + a_j \frac{\partial \sigma}{\partial a_j} = v(j)^2 \sigma. \quad (6)$$

Now set  $j = 1, 2, \dots, n$ , in turn, to find

$$\begin{aligned} a_1 \frac{\partial \sigma}{\partial a_1} &= \sigma v(1)^2, \\ b_1 \frac{\partial \sigma}{\partial b_1} &= \sigma (u(1)^2 - v(1)^2), \\ a_k \frac{\partial \sigma}{\partial a_k} &= \sigma \left[ \sum_{i=1}^k v(i)^2 - \sum_{j=1}^{k-1} u(j)^2 \right], \\ b_k \frac{\partial \sigma}{\partial b_k} &= \sigma \sum_{i=1}^k (u(i)^2 - v(i)^2). \end{aligned}$$

Since  $\mathbf{u}^t \mathbf{u} = \mathbf{v}^t \mathbf{v} = 1$ , the complementary expressions follow from  $\sum_{i=1}^k v(i)^2 = 1 - \sum_{j=k+1}^n v(j)^2$ . In particular,  $a_n \partial \sigma / \partial a_n = \sigma u(n)^2$ .  $\square$

**Corollary 1.** For  $k = 1, 2, \dots, n$ ,

$$\begin{aligned} \left| \frac{\partial \sigma}{\partial a_k} \cdot \frac{a_k}{\sigma} \right| &< 1, \quad \left| \frac{\partial \sigma}{\partial b_k} \cdot \frac{b_k}{\sigma} \right| < 1, \\ 0 &\leq \frac{\partial \sigma}{\partial a_k} \cdot \frac{a_k}{\sigma} + \frac{\partial \sigma}{\partial b_k} \cdot \frac{b_k}{\sigma} = u(k)^2 < 1. \end{aligned}$$

**Proof.** Expressions of the form  $\sum v(i)^2 - \sum u(j)^2$  must lie in  $[-1, +1]$ . The assumption that  $a_k b_k \neq 0$ ,  $k = 1, \dots, n-1$ , ensures that equality is never obtained, e. g.  $|v(1)| < 1$ .  $\square$

The formulae in Theorem 1 show that we may have

$$\frac{\partial \sigma}{\partial a_k} \cdot \frac{a_k}{\sigma} \ll 1, \quad \frac{\partial \sigma}{\partial b_k} \cdot \frac{b_k}{\sigma} \ll 1.$$

Examples of very small relative condition numbers are given in Section 6 and show that this phenomenon is common and is not confined to the positive definite case.

**Corollary 2.** The eigenvalues of a symmetric tridiagonal matrix that is positive definite are determined to high relative accuracy by the entries in the Cholesky factor.

#### 4. A classical example

A system of masses connected by springs is a standard example in courses on dynamics. The equations of motion  $M\ddot{x} + Kx = \mathbf{0}$  give rise to an eigenvalue problem  $(K - \omega^2 M)v = \mathbf{0}$  with  $M = \text{diag}(m_1, \dots, m_n)$ ,  $x = v \exp(i\omega t)$  and

$$K = \text{tridiag} \begin{pmatrix} & -k_2 & -k_3 & \cdot & -k_n \\ k_1 + k_2 & & k_2 + k_3 & \cdot & \\ & -k_2 & -k_3 & \cdot & -k_n \\ & & & & k_n \end{pmatrix}$$

where the  $k_i$  are the spring constants. The point is that  $K$  may be factored as

$$K = BK_d B^t,$$

where  $K_d = \text{diag}(k_1, \dots, k_n)$  and

$$B = \text{bidiag} \begin{pmatrix} & -1 & -1 & \cdot & -1 & -1 \\ 1 & & 1 & \cdot & & \\ & & & & 1 & 1 \end{pmatrix}.$$

It follows that the desired frequencies  $\omega$  satisfy

$$\left[ \left( M^{-1/2} B K_d^{1/2} \right) \left( M^{-1/2} B K_d^{1/2} \right)^t - \omega^2 I \right] u = \mathbf{0}.$$

So the  $\omega$  are the singular values of  $U = M^{-1/2} B K_d^{1/2}$  and each  $\omega$ , however small, is determined to high relative accuracy by the entries of  $U$ . In contrast the small  $\omega^2$  are not always determined to high relative accuracy by the entries of  $M^{-1/2} K M^{-1/2}$  where the additions on the diagonal discard important information. The tactical error, passed down from teacher to student for decades, was to use  $K$  instead of  $B K_d B^t$ , the product form. There are no adds or subtracts in forming  $U$  whose  $j$ th row contains  $a_j = (k_j/m_j)^{1/2}$ ,  $b_j = -(k_{j+1}/m_j)^{1/2}$ . See [12]. Current research considers whether the  $B$ ,  $K_d$ , and  $M$  matrices that arise in 2D and 3D applications also define the small eigenvalues to high relative accuracy.

#### 5. Symmetric tridiagonals

If  $T$  is not positive definite (p.d.) but  $T + \sigma I$  is p.d. then Theorem 1 permits us to say that the quantities  $(\lambda + \sigma)$  are determined to high relative accuracy by the Cholesky factor of  $T + \sigma I$ . This is not very satisfactory and so the next step is to see to what extent Theorem 1 can be extended.

All tridiagonals  $T$  may be written as

$$T = L\Omega L^t, \tag{7}$$

where  $L$  is lower triangular with positive diagonal and  $\Omega$  is a signed, symmetric permutation matrix.  $\Omega$  will be a direct sum of  $\pm 1$  and

$$\pm \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

entries. Note that  $\Omega^2 = I$ . The spectral factorization of invertible  $T$  is written as

$$T = SAS^t, \quad I = SIS^t \quad (8)$$

and each eigenvalue  $\lambda$  is written as

$$\lambda = \text{sign}(\lambda)\sigma^2, \quad (9)$$

so that  $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_n)$  is positive definite. Note that Eq. (8) is unique only up to permutations of  $A$ , i.e. the labelling of eigenvalues.

Our interest is in how the eigenvalues of  $T = L\Omega L^t$  respond to small relative changes to the nonzero entries of  $L$  while  $\Omega$  remains fixed. Theorem 1 used the singular vectors of  $L$  to express the sensitivity when  $\Omega = I$ . When  $\Omega \neq I$  it is useful to introduce an indefinite, or improper, SVD of  $L$ , namely

$$L = S\Sigma P^t, \quad (10)$$

where  $S$  and  $\Sigma$  are given by Eqs. (8) and (9) and  $P$  is  $\Omega$ -sign( $A$ )-orthogonal, i.e.

$$P^t\Omega P = \text{sign}(A). \quad (11)$$

When some  $\Omega \neq I$  then  $\Omega$  defines an indefinite ‘inner product’  $\langle u, v \rangle := v^t\Omega u$  that satisfies all the axioms except for positivity.

From Eq. (10),

$$P = L^t S \Sigma^{-1}, \quad (12)$$

we find

$$\begin{aligned} P^t\Omega P &= \Sigma^{-1} S^t L \Omega L^t S \Sigma^{-1} \\ &= \Sigma^{-1} S^t T S \Sigma^{-1} \\ &= \text{sign}(A), \quad \text{by Eq. (9),} \end{aligned}$$

verifying Eq. (11).

**Theorem 2.** *An invertible symmetric tridiagonal matrix  $T$  may be written as  $T = L\Omega L^t$  with*

$$L = \text{bidiag} \begin{pmatrix} a_1 & a_2 & \cdot & a_{n-1} & a_n \\ & b_1 & \cdot & \cdot & b_{n-1} \end{pmatrix}.$$

$\Omega = a$  signed symmetric perturbation matrix.

Let the spectral factorization of  $T$  be given by Eqs. (8) and (9). Let the columns of  $P = (p_1, \dots, p_n)$ , defined by Eq. (12), be the right  $\Omega$ -singular vectors of  $L$  and note that the columns of  $S$  are the left singular vectors of  $L$ .

Let  $(\sigma, s^t, \Omega p)$  be a typical improper singular triple of  $L$  and let the corresponding eigenvalue of  $T$  be  $\text{sign}(\lambda)\sigma^2$ . Let  $\Omega p(k)$  denote the  $k$ th entry of  $\Omega p$ , then

$$\begin{aligned}\frac{a_1}{\sigma} \cdot \frac{\partial \sigma}{\partial a_1} &= s(1)^2, \\ \frac{b_1}{\sigma} \cdot \frac{\partial \sigma}{\partial b_1} &= \text{sign}(\lambda)p(1)\Omega p(1) - s(1)^2, \\ \frac{a_k}{\sigma} \cdot \frac{\partial \sigma}{\partial a_k} &= \sum_{i=1}^k s(i)^2 - \text{sign}(\lambda) \sum_{j=1}^{k-1} p(j)\Omega p(j) \\ &= \text{sign}(\lambda) \sum_{i=k}^n p(i)\Omega p(i) - \sum_{j=k+1}^n s(j)^2, \\ \frac{b_k}{\sigma} \cdot \frac{\partial \sigma}{\partial b_k} &= \text{sign}(\lambda) \sum_{i=1}^k p(i)\Omega p(i) - \sum_{j=1}^k s(j)^2 \\ &= \sum_{i=k+1}^n s(i)^2 - \text{sign}(\lambda) \sum_{j=k+1}^n p(j)\Omega p(j), \\ \frac{a_n}{\sigma} \cdot \frac{\partial \sigma}{\partial a_n} &= \text{sign}(\lambda)p(n)\Omega p(n).\end{aligned}$$

**Proof.** Multiply each side of  $L = S\Sigma P^t$  by  $\Omega P$  and use Eq. (11) to find  $L\Omega P = S\Sigma \text{sign}(\Lambda)$ , or

$$L\Omega p_i = s_i \sigma_i \text{sign}(\lambda_i), \quad i = 1, 2, \dots, n. \quad (13)$$

From the definition of  $L$

$$L^t s_j = P\Sigma S^t s_j = P\Sigma e_j = p_j \sigma_j, \quad j = 1, 2, \dots, n. \quad (14)$$

By Eq. (11),

$$p_j^t \Omega p_j = \text{sign}(\lambda_j), \quad s_j^t s_j = 1. \quad (15)$$

Premultiply Eq. (13) by  $s_i^t$ , write out the expression  $\text{sign}(\lambda_i)\sigma_i = s_i^t L\Omega p_i$  to see that it is linear in  $a_j$  and  $b_j$ . Drop the index  $i$  to find

$$\text{sign}(\lambda) \frac{\partial \sigma}{\partial a_j} = \Omega p(j)s(j), \quad \text{sign}(\lambda) \frac{\partial \sigma}{\partial b_j} = \Omega p(j)s(j+1). \quad (16)$$

Now use the bidiagonal form of Eqs. (13) and (14). For  $b_0 = b_n = 0$  and  $j = 1, 2, \dots, n$ ,

$$\begin{aligned}b_{j-1}\Omega p(j-1) + a_j\Omega p(j) &= \text{sign}(\lambda)s(j)\sigma, \\ a_j s(j) + b_j s(j+1) &= p(j)\sigma.\end{aligned}$$

Multiply by  $s(j)$  and  $\Omega p(j)$ , respectively, and use Eq. (16) to obtain



$$\begin{aligned}\operatorname{sign}(\lambda)b_{j-1}\frac{\partial\sigma}{\partial b_{j-1}} + \operatorname{sign}(\lambda)a_j\frac{\partial\sigma}{\partial a_j} &= \operatorname{sign}(\lambda)s(j)^2\sigma, \\ \operatorname{sign}(\lambda)a_j\frac{\partial\sigma}{\partial a_j} + \operatorname{sign}(\lambda)b_j\frac{\partial\sigma}{\partial b_j} &= p(j)\Omega p(j)\sigma.\end{aligned}$$

Now set  $j = 1, 2, \dots, n$  in turn to obtain the conclusion of Theorem 2.  $\square$

**Remark 1.** From Eq. (9), it follows that, for  $\lambda \neq 0$ ,

$$\frac{a_j}{\lambda} \frac{\partial \lambda}{\partial a_j} = 2 \frac{a_j}{\sigma} \frac{\partial \sigma}{\partial a_j}, \quad \frac{b_j}{\lambda} \frac{\partial \lambda}{\partial b_j} = 2 \frac{b_j}{\sigma} \frac{\partial \sigma}{\partial b_j}.$$

**Remark 2.** From Eq. (12),  $\mathbf{p} = L^1 \mathbf{s} / \sigma$  and  $\|\mathbf{p}\| \leq \|L\| / \sigma$  but this inequality is rather crude. Since  $L\Omega L^1 \mathbf{s} = \mathbf{s}\lambda$  an alternative expression for  $\mathbf{p}$  is

$$\mathbf{p} = \Omega L^{-1} \mathbf{s} \cdot \sigma \cdot \operatorname{sign}(\lambda)$$

and so

$$\|\mathbf{p}\| \leq \sigma \|L^{-1}\|.$$

Thus

$$\|\mathbf{p}\| \leq \sqrt{\operatorname{cond}(L)} = \sqrt{\|L\| \cdot \|L^{-1}\|}.$$

But even this bound ignores cancellation in  $L^1 \mathbf{s}$  and  $L^{-1} \mathbf{s}$ .

**Remark 3.** Section 4 indicates that there are important applications where the matrix of interest is given in product form and the small eigenvalues may be defined to high relative accuracy by the factors but not by their product. When a given product does not yield the desired accuracy it is recommended to look for better product representations and that is focus of our current investigations.

The message of Theorem 2 is that  $L$  certainly determines an eigenvalue  $\lambda$  of  $L\Omega L^1$  to high relative accuracy when the associated  $\mathbf{p}$ -vector has entries bounded by a modest value such as 10.

The point to notice is that when  $\Omega = I$  the vector  $\mathbf{p} = \Omega L^{-1} \mathbf{s} \sigma$  satisfies  $\|\mathbf{p}\| = 1$  however ill-conditioned  $L$  may be. Thus we may expect that, even when  $\Omega \neq I$ , the bound  $\|\mathbf{p}\| \leq \sqrt{\operatorname{cond}(L)}$  is far from tight in many, but not all, cases. The next section supplies some evidence.

## 6. Examples

Description of the tables given below. For each  $\lambda \neq 0$  of each  $n \times n$  tridiagonal  $T$  there are  $2n - 1$  relative condition numbers ( $\operatorname{rcond}$ ),  $|(a_j/\lambda)(\partial\lambda/\partial a_j)|$ ,



(2b)  $T$  is the  $[1,2,1]$  matrix shifted by  $(\text{eig}(20)+\text{eig}(21))/2 = 1.9252$ .

$$n = 41, \quad \|L\| = 7.207, \quad \text{cond}(L) = 335.5,$$

[illegible]

$\lambda$	$\ p\ $	max rcond	w.r.t	min rcond	w.r.t.
0.0747	8.74	33.2	41	$3.22 \times 10^{-15}$	4
-0.0747	8.74	33.0	-40	$2.72 \times 10^{-14}$	21
0.0747	8.74	33.2	41	$3.22 \times 10^{-15}$	4
2.07	1.11	0.0994	-17	0.000474	40

(2c)  $T$  is the  $[1,2,1]$  matrix shifted by  $\text{eig}(5) = 0.13825$

$$n = 41, \quad \|L\| = 3.164, \quad \text{cond}(L) = 2.336 \times 10^8,$$

[illegible]

$\lambda$	$\ P\ $	max rcond	w.r.t	min rcond	w.r.t
$1.65 \times 10^{-16}$	1.05	2.22	41	0.01272	-1
-0.0494	5.07	8.38	-16	$4.44 \times 10^{-15}$	41
-0.0494	5.07	8.38	-16	$4.44 \times 10^{-15}$	41
3.72	1.01	0.0494	29	$9.86 \times 10^{-17}$	41

(3a) The tridiagonal is  $W_{21}^+$ . See [13], p. 330.

$$n = 21, \quad \|L\| = 3.278, \quad \text{cond}(L) = 10.49,$$

$$\Omega = \text{diag}(+++++ - +++++).$$

$\lambda$	$\ p\ $	max rcond	w.r.t	min rcond	w.r.t.
0.254	1.39	1.528	12	$3.62 \times 10^{-14}$	11
0.948	1.35	1.470	-10	$6.34 \times 10^{-12}$	-1
-1.13	1.72	3.26	11	$1.034 \times 10^{-15}$	1
3.04	1.01	0.352	7	$2.24 \times 10^{-8}$	-1

This well-known matrix has one negative eigenvalue. It is the most sensitive but nevertheless is extremely well determined as is guaranteed by the small value of  $\text{cond}(L)$ .

(3b)  $W_{21}^+$  shifted by  $\text{eig}(18) = 9.210678647304919$ .

$$n = 21, \quad \|L\| = 1810, \quad \text{cond}(L) = 2.376 \times 10^{11},$$

$$\Omega = \text{diag}(+ - - - - - - - - - + - - - - - - - - + -).$$

$\lambda$	$\ p\ $	max rcond	w.r.t	min rcond	w.r.t.
$-2.8 \times 10^{-16}$	2.61	5.28	11	0.292	1
$5.64 \times 10^{-11}$	605	$3.66 \times 10^5$	11	$5.18 \times 10^{-6}$	21
$5.64 \times 10^{-11}$	605	$3.66 \times 10^5$	11	$5.18 \times 10^{-6}$	21
-7.42	1	0.354	14	$3.4 \times 10^{-26}$	21

(3c)  $W_{21}^+$  shifted by  $\text{eig}(19) = 9.210678647361332$ .

$$n = 21, \quad \|L\| = 3.215, \quad \text{cond}(L) = 2.405 \times 10^8$$

$$\Omega = \text{diag}(+ - - - - - - - - - - - - - - - - - - + -).$$

$\lambda$	$\ p\ $	max rcond	w.r.t	min rcond	w.r.t.
$6.44 \times 10^{-16}$	0.527	1.708	20	0.292	-1
$-5.64 \times 10^{-11}$	1	1.0	7	$6.34 \times 10^{-6}$	21
$6.44 \times 10^{-16}$	0.527	1.708	20	0.292	-1
-6.17	1	0.348	13	$6.38 \times 10^{-24}$	21

(3b) and (3c) show the extreme sensitivity to a shift. The shifts agree to 11 decimals, both matrices are almost singular and yet (3c) yields very small condition numbers while (3b) yields very large ones. We have no explanation of this phenomenon.

## 7. Nonsymmetric tridiagonals

Any real tridiagonal is diagonally similar to  $\Delta T$ , where  $T$  is real symmetric and  $\Delta = \text{diag}(\delta_1, \dots, \delta_n)$ ,  $\delta_i = \pm 1$ . This similarity transformation is an instance of balancing a nonsymmetric matrix; a simple and often used way to reduce the norm of the original matrix. We assume that the eigenvalues  $\lambda_i$  are distinct but we allow them to be complex. Thus

$$Ts_i = As_i\lambda_i, \quad i = 1, \dots, n, \quad (17)$$

$$S = [s_1, \dots, s_n], \quad \text{possibly complex.} \quad (18)$$

Any such  $T$  may be written as

$$T = L\Omega L^t, \quad (19)$$

where  $L$  is lower triangular and  $\Omega$  is a sign symmetric permutation (s.s.p.) matrix. For simplicity we shall assume that  $\Omega$  is diagonal. A convenient normalization for the eigenvectors is

$$S^t \Delta S = I. \quad (20)$$

Let  $A = \text{diag}(\lambda_1, \dots, \lambda_n)$  and denote by  $A^{1/2}$  the principal square root of  $A$ ;  $\lambda = \rho \exp(2i\theta)$ ,  $\lambda^{1/2} = \rho \exp(i\theta)$ ,  $-\pi/2 < \theta \leq \pi/2$ . As in the symmetric case there is an improper singular value decomposition of  $L$ ;

$$L = \Delta S A^{1/2} P^t, \quad (21)$$

defining  $P$  which is complex in general. Note that

$$\begin{aligned} P^t \Omega P &= A^{-1/2} (\Delta S)^{-1} L \Omega L^t (\Delta S)^{-1} A^{-1/2}, \quad \text{by Eq. (21),} \\ &= A^{-1/2} (\Delta S)^{-1} T (\Delta S)^{-t} A^{-1/2}, \quad \text{by Eq. (19),} \\ &= A^{-1/2} [\Delta S^{-1} (\Delta S)^{-t}] A^{-1/2}, \quad \text{by Eq. (17) = } I, \text{ by Eq. (20).} \end{aligned} \quad (22)$$

With  $\Delta$  and  $\Omega$  fixed we study how  $A$  depends on  $L$ . Write

$$L = \text{bidiag} \begin{pmatrix} a_1 & & & a_{n-1} & a_n \\ & b_1 & & & \\ & & b_2 & & \\ & & & \ddots & \\ & & & & b_{n-1} \end{pmatrix}$$

and consider a typical ‘singular’ triple  $(\lambda, s, p^t)$ ;

$$\begin{aligned} L(\Omega p) &= (\Delta s) \lambda^{1/2}, \quad (23) \\ s^t L &= \lambda^{1/2} p^t. \end{aligned} \quad (24)$$

Now we can state the relative sensitivity of the eigenvalues of  $\Delta L \Omega L^t$  to  $L$ ’s entries. The expressions may be complex.

**Theorem 3.** *If  $\Delta T$  has distinct eigenvalues, and Eq. (19) holds with  $\Omega = \text{diag}(\omega_1, \dots, \omega_n)$ ,  $\omega_i = \pm 1$ , then, for  $\lambda \neq 0$ ,*

$$\frac{1}{2} \frac{a_j}{\lambda} \frac{\partial \lambda}{\partial a_j} = \sum_{k=1}^j \delta_k s(k)^2 - \sum_{m=1}^{j-1} \omega_m p(m)^2 = \sum_{m=j}^n \omega_m p(m)^2 - \sum_{k=j+1}^n \delta_k s(k)^2,$$

$$\frac{1}{2} \frac{b_j}{\lambda} \frac{\partial \lambda}{\partial b_j} = \sum_{m=1}^j [\omega_m p(m)^2 - \delta_m s(m)^2] = \sum_{m=j+1}^n [\delta_m s(m)^2 - \omega_m p(m)^2].$$

**Proof.** Premultiply Eq. (23) by  $s^t$  and use Eq. (20) to find

$$\lambda^{1/2} = s^t L \Omega p.$$

Thus  $\lambda^{1/2}$  depends linearly on  $a_j$  and  $b_j$ , the entries of  $L$ , and

$$\frac{\partial}{\partial a_j} (\lambda^{1/2}) = s(j) \omega_j p(j), \quad (25)$$

$$\frac{\partial}{\partial b_j} (\lambda^{1/2}) = s(j+1) \omega_j p(j). \quad (26)$$

Next write out Eqs. (23) and (24) in detail,

$$b_{j-1} \omega_{j-1} p(j-1) + a_j \omega_j p(j) = \delta_j s(j) \lambda^{1/2}, \quad s(j) a_j + s(j+1) b_j.$$

Multiply the first equation by  $s(j)$  and the second one by  $\omega_j p(j)$  and substitute the expressions from Eqs. (25) and (26) to find for  $j = 1, 2, \dots, n$  and  $b_0 = b_n = 0$ ,

$$b_{j-1} \frac{\partial}{\partial b_{j-1}} (\lambda^{1/2}) + a_j \frac{\partial}{\partial a_j} (\lambda^{1/2}) = \delta_j s(j) \lambda^{1/2}, \quad (27)$$

$$a_j \frac{\partial}{\partial a_j} (\lambda^{1/2}) + b_j \frac{\partial}{\partial b_j} (\lambda^{1/2}) = \lambda^{1/2} \omega_j p(j)^2. \quad (28)$$

Finally observe that

$$\left( \frac{a_j}{\lambda^{1/2}} \right) \frac{\partial}{\partial a_j} (\lambda^{1/2}) = \frac{1}{2} \frac{a_j}{\lambda} \frac{\partial \lambda}{\partial a_j}.$$

The system of equations (27) and (28) is triangular and may be solved by setting  $j = 1, 2, \dots$ , in Eqs. (27) and (28) to give the first expressions for the sensitivities. Solving in reverse order gives the second expressions.  $\square$

Theorem 3 shows that an eigenvalue (possible complex) of the pair  $L \Omega L^t$ ,  $\lambda$  is defined to high relative accuracy whenever its singular vectors  $p$  and  $s$  have modest norms, say  $\|p\| < 10$ ,  $\|s\| < 10$ . Note that Eq. (20) does not guarantee that  $\|s\| = 1$ .

## Acknowledgements

I would like to thank Inderjit Dhillon for producing the numerical examples in Section 6.

## References

- [1] B.N. Parlett, *The Symmetric Eigenvalue Problem*, 2nd edition, SIAM, Philadelphia, 1997.
- [2] P. Deift, J.W. Demmel, L.-C. Li, C. Tomei, The bidiagonal singular value decomposition and Hamiltonian mechanics, *SIAM J. Numer. Anal.* 28 (1991) 1463–1516.
- [3] J.W. Demmel, K. Veselić, Jacobi's method is more accurate than QR, *SIAM J. Matrix Anal. Appl.* 13 (1992) 1204–1245.
- [4] I. Slapnicar, *Accurate Symmetric Eigenreduction by a Jacobi Method*, Thesis, Fernuniversitat Hagen, Germany, 1992.
- [5] K. Veselić, I. Slapnicar, Floating point perturbations of Hermitian matrices, *Linear Algebra Appl.* 195 (1993) 81–116.
- [6] J.W. Demmel, W. Kahan, Accurate singular values of bidiagonal matrices, *SIAM J. Sci. Stat. Comput.* 11 (1990) 873–912.
- [7] J. Barlow, J.W. Demmel, Computing accurate eigensystems of scaled diagonally dominant matrices, *SIAM J. Numer. Anal.* 27 (1990) 762–791.
- [8] S.C. Eisenstat, I.C.F. Ipsen, Relative perturbation techniques for singular value problems, *SIAM J. Numer. Anal.* 32 (1995) 1972–1988.
- [9] V. Hari, Z. Drmač, On scaled almost diagonal Hermitian matrix pairs, *SIAM J. Matrix. Anal. Appl.*, to appear.
- [10] R.-C. Li, Relative perturbation theory: (I) eigenvalue and singular value variations, Tech. Rep. Univ. of California, Berkeley, UCB/CSD-94-855, Computer Science Division, Department of EECS, Univ. of California, Berkeley, 1994. (Revised January 1996). to appear *SIAM Matrix Anal. Appl.*
- [11] J.W. Demmel, *Applied Linear Algebra*, SIAM, Philadelphia, 1997.
- [12] K. Veselić, On linear vibrational systems with one-dimensional damping, *Appl. Anal.* 29 (1988) 1–18.
- [13] J.H. Wilkinson, *The Algebraic Eigenvalue Problem*, Clarendon Press, Oxford, 1965.